

統計数値データの圧縮と 国勢調査データベース・システム

内 藤 三 義

はじめに

本稿ではコンピュータで扱う統計数値データを対象にした、データ圧縮の新しい技法を説明し、併せて、国勢調査データファイルの総合的データベース化の方法を提案する。

筆者は既に、「日経財務データ」、「家計調査年報」などを対象にした、数値データ圧縮法について私見を述べたが¹⁾、特定の種類の統計データだけでなく、たいいていの統計数値に「汎用」的に利用できるとされるデータ圧縮法へと、これまでの圧縮法を「改良」した。この新しい「圧縮技法」を説明するのが、本稿の第1の目的である。

第2に、国勢調査データについてであるが、磁気テープ媒体でも提供されているこのデータ

をコンピュータで直接的に利用することは、これまでは利用環境の制約などから、研究者の間にあまり普及してはいなかったようである。このデータを利用する機会は、社会科学分野では非常に多いが、購入費用、データベース・システムとする時に必要な機器施設環境、検索システム作成などの諸問題がネックとなり、一般の教育・研究での利用がなかなか進まないのが現状である。そこで、利用したい磁気テープ媒体などのデータがあれば、比較的簡単にそれをパーソナルにも処理できるデータベース・システムとして作成、稼働させるソフトを開発した。このソフト・システムの内容を示すのが、本稿の第2の目的である。

1. 統計数値圧縮法

コンピュータで扱うデータが大容量化するにつれ、より効率的なデータ圧縮技法がますます求められているのは周知のことである。加速度的に進行しているコンピュータ関連技術の革新によって、コンピュータの処理速度の向上、記憶装置のコンパクト化と大容量化、通信速度の高速化が進められ、ますます大きな情報をより速く処理できるようになってきているが、それらをより有効に稼働させるためにも、より正確で効率的なデータ圧縮が求められている。

「データ圧縮技法は本来情報に付帯する冗長性を取り除いて、最大限に環境条件に適したデータの格納・伝送を行うための技術」²⁾であり、アルファベット文字を表すための平均コード長を圧縮する「ハフマンのコード化」³⁾など多様な技法⁴⁾がすでに開発されている。しかし

それらは文書やプログラムなどファイル般を対象にしたものであるもので、それらだけを利用することでは、統計数値データに対する高い圧縮を実現することはできない。数値の表現方法としては、バイナリー型で整数や浮動小数点数値を表現する方式や、10進数の1桁を0.5バイトで表現するパック化という表現方法あり、これらは10進数1文字・1バイトで表現するキャラクター型の表現からみれば1つの「圧縮」方法であるとも言える。しかし、これらも、数値の性質が一定している場合に、キャラクター型に対しておおよそ50%程度まで圧縮出来る程度で、高率の圧縮を統計数値一般に保証できるものではない。そこで、統計数値(統計書に記載されるような数値)に限定しながら効果的な圧縮法を考えてみよう。

1-1. 統計数値の特性

統計書に記載され、磁気テープなどの記憶媒体に保存される数値の特徴としては次の5つのことをあげることができる。

第1には、数値は概ね10進数で表記される正負実数であるが、整数だけでなく小数点以下が表示されている数値もあるということである。磁気テープなどへの記録のされ方は、パック化されているものもあるが、小数点以下の桁数表示がランダムに変化するケースが含まれている統計書もあり、その場合データ形式の汎用性を確保するためにもEBCDICコードによるピクチャー形式で記録されているものが多い。

第2に10進数等で表示された数値の有効桁数あるいは最大桁数であるが、これは統計書によって一定していないというものの、9桁を最大（マイナス記号を1桁とすると10桁）とするものが非常に多い。日本の場合、総人口が「億」という10進数9桁の値であることを考えてみれば、概ね9桁以下が多いのは想定もつく。但しそれ以上の最大桁数のものも無いわけでない。

第3には、数値は統計書では表の形で記載されることが多く、磁気テープではその表の1行、もしくは1列が1レコードという単位にまとめられることが多い。その場合レコードの意味を示すためのRID (Record Identifier) がレコードの先頭に固定書式で記録されることが多い。

第4には、該当数値データがないとか該当数値が計算できない、あるいは極めて小さいか大きいといった場合に、何らかの欠損値表示が統計書では行われ、磁気テープ・データでも特殊な文字列やコードがこれを示すために使われる。この「欠損値」の全データに占める比率は決して低いものではない。

第5は数値の特性というより、データファイルの特性とも言えるが、この種の数値データは、静的で書き直しのされないデータであり（統計の間違いによって、後に「変更」される可能性が無いわけではないが）、別の言葉で言

えば読みだし専用、もしくは追記型のデータという性格を持っている。一時的な保存のためだけでなく、恒常的に利用するデータベース・ファイルを圧縮しても問題が起きないのはこの特徴の故である。

以上のような諸特徴が統計数値データには考えられるが、これらの中でデータ圧縮のために考えなければならない問題は、(1)正負数値の混在、(2)整数と小数点表示の値の混在、(3)有効最大桁数が概ね10進数9桁であること、(4)かなりの出現頻度ある欠損値の存在、(5)レコードを単位とした数値データの記録（数値が1項目つづでなく、何項目かまとめて記録されている）、という5つの問題であり、データベース・システムとして構成するときには、(6)RID情報をどう効率的に整理するかという問題が追加される。

では、実際の数値はどのような値の分布になっているのだろうか。家計調査年報・昭和63年、昭和60年国勢調査基本集計・京都府、日経財務データ（1988年まで）を例に、その数値の精度を落とすこと無く表現出来るビット列数と、レコードを単位として見た正負値の混在数、小数点表示の数値の割合と、その桁数などについて見てみると表1のようになった。

3つのデータの、有効最大桁数は10進数9桁であり、実際に記録されている数値は29 bit以下で表現出来るが、15 bit以下でも表現出来る数値が全体の90%を越えていることが分かる。レコードを単位としてみると、全て正であるものの、整数を含め繰下がり桁数がレコード中の数値で一定であるものがかなりの率を占めることが分かる。

1-2. N式数値圧縮法

そこでデータ圧縮のために、数値の表現方法を次のように整理しよう。

(1) 正負の判定はレコード単位、もしくは個別数値別に情報を与える。

(2) 小数点付の数値（10進数）はそれを小数点以下の有効桁数分繰り上げて整数化し、その繰上げ桁数を判定できる情報を付加する⁵⁾。

表－1 諸統計データで現れる数値

(1) 数値を表現するために必要なビット数

	家計調査年報		国勢調査		日経財務	
数値総数	1245619	100.000	810820	100.000	14388465	100.000
欠損値	247917	19.903	180392	22.248	6424285	44.649
0 bit	31661	2.542	-	-	178419	1.240
1 bit	2951	0.237	33893	4.180	159456	1.108
2 bit	5444	0.437	35117	4.331	264012	1.835
3 bit	9012	0.723	37167	4.584	374479	2.603
4 bit	11229	0.901	40252	4.964	472326	3.283
5 bit	15557	1.249	45922	5.664	602439	4.187
6 bit	20773	1.668	53733	6.627	637024	4.427
7 bit	29009	2.329	63348	7.813	703765	4.891
8 bit	48671	3.907	59824	7.378	777182	5.401
9 bit	73786	5.924	64498	7.955	747923	5.198
10 bit	104640	8.401	48334	5.961	692167	4.811
11 bit	135994	10.918	38014	4.688	615712	4.279
12 bit	134748	10.818	34228	4.221	519061	3.607
13 bit	123344	9.902	26112	3.220	403499	2.804
14 bit	95334	7.654	19259	2.375	289394	2.011
15 bit	54739	4.395	12935	1.595	185844	1.292
16 bit	30762	2.470	8107	1.000	111277	0.773
17 bit	20883	1.677	5330	0.657	61603	0.428
18 bit	15359	1.233	2232	0.275	30725	0.214
19 bit	16925	1.359	1146	0.141	32210	0.224
20 bit	7324	0.588	604	0.074	6509	0.045
21 bit	3833	0.308	297	0.037	2349	0.016
22 bit	2344	0.188	76	0.009	575	0.004
23 bit	1701	0.137	-	-	148935	1.035
24 bit	1150	0.092	-	-	-	-
25 bit	164	0.013	-	-	-	-
26 bit	134	0.011	-	-	-	-
27 bit	204	0.016	-	-	-	-
28 bit	26	0.002	-	-	-	-
29 bit	1	0.000	-	-	-	-

(2) レコードを単位としてみた数値の性質

	家計調査年報		国勢調査		日経財務	
レコード総数	5651	100.000	96416	100.000	52705	100.000
正・整数	1211	21.430	82662	85.735	2709	5.140
正・小数点1桁	4	0.071	2353	2.440	-	-
正・小数点2桁	379	6.707	2762	2.865	-	-
正・小数点3桁	58	1.026	-	-	-	-
正・小数点0～1桁	45	0.796	991	1.028	-	-
正・小数点0～3桁	2484	43.957	7648	7.932	-	-
正負混在・整数	-	-	-	-	49996	94.860
混在・小数点1桁	29	0.513	-	-	-	-
混在・小数点0～1桁	391	6.919	-	-	-	-
混在・小数点0～3桁	1050	18.581	-	-	-	-

(3) 整数化された全ての数値はその絶対値を2進数のビット列で表現するが、固定ではなく、可変ビット列で表し、そのビット列数を判定できる情報を付加する、なお、0(ゼロ)以外

を表現するビット列の先頭には必ず1が表れるがこの1ビットは省略する。

(4) 欠損値を表現する特殊なコードを設定する。

1-2-1. N式整数の構造

上記のことを実現するような数値の表現方法として、具体的には次のような「N式整数」を規定する。

数値を表すビット列の並びを、 $[I_1 \sim I_i K_1 \sim K_k M D_1 \sim D_d]$ とし、その意味を次のように設定する。

- $I_1 \sim I_i$: 欠損値、もしくは数値を表現するビット数 (d の値) を示す
 $K_1 \sim K_k$: 10進数換算の繰り上げ桁数を示す (下に説明するレコード単位制御情報RCI-Nを付加しないとき、このビットは通常3ビットとする)
 M : 正/負を示す
 $D_1 \sim D_d$: 先頭に $[1] b$ を付加したビット列が正整数化された数値の値である

ここで $I_1 \sim I_i$ のビット列の意味は次のようにする。

- 11 : 欠損値を表す。(KMD) のビット列は与えられない
0 xxxx : $d = [xxxx] b - 2$
但 $[xxxx] b = 0$ はゼロ
10xxxx : $d = 14 + [xxxx] b$
但 $[xxxx] = [1111]$ の時は更に I 情報が 4 bit 付加され、 $d = [xxxx] b + 29$ とする。

1-2-2. N式レコード制御情報 (RCI-N0とRIC-N1の構造)

さらに、レコードを単位とした正負の混在や繰上げ桁数を示すために、次のような「N式レコード制御情報 (RCI-N0) (1 or 2) + 5 ビットを、レコードの先頭に置く場合もあるとする。

(1) RCI-N0

(1)-1 正負制御情報 (1 or 2) ビット

- 0 : レコード内の数値は全て正
10 : レコード内の数値は全て負
11 : 正負混在のため、0 と欠損値以外では正負制御ビットとして各数値データ内に1ビットを確保している。

(1)-2 繰上げ桁数情報 5 ビット

- 0 xxxx : レコード内の繰上げ桁数は全て同じで $[xxxx] b$ 桁
1 xxxx : レコード内の繰上げ桁数は全て同じではない。欠損値以外では桁数情報を与えるため $[xxxx] b$ ビットを各数値毎に確保している。

(2) RCI-N1

なお、この情報は、次のような「N式レコード制御変形情報 (RCI-N1)」として書くこともあるとする。

- 00 : レコード内の数値は全て正整数
01 : レコード内の数値は全て正で、繰上げ桁数情報のため 2 bit を確保している
100 : レコード内の数値は全て正で、全て10進数2桁繰上げて入力されている
101 : レコード内の数値は全て正で、繰上げ桁数情報のため 1 bit を確保している
110 : レコード内の数値は全て正で、全て10進数1桁繰上げて入力されている
111 : RCI-N0の情報がこの3ビットの後に与えられている

1-2-3. N式数値圧縮法表現の具体例

前述のRCI-N1付のN式整数が一般的なN式数値圧縮法であるが、この方法によって欠損値をふくめた9つの数値からなる1レコードの数値データを表現すると、表2のようになる。事例2の場合、9つの数値を表現する1レコード12バイトは、16進数表記すると $[98\ 02\ 79\ 11\ 21\ 25\ 89\ 90\ 52\ F1\ 60\ 08] h$ となる。

このN式数値圧縮法を利用して、昭和60年国勢調査基本集計・京都府の圧縮効率を見たところ

表2 ある1レコードの数値の表現

事 例 1					事 例 2				
RCI-N1	111 11 10010				RCI-N1	100			
数値	I	K	M	D	数値	I	K	M	D
欠損値	11...	欠損値	11...
0	00000	00	0	0.00	00000
1	00001	00	0	0.01	00001
2	00010	00	0	0.....	1.00	00111	..	.	100100....
3	00010	00	0	1.....	2.00	01000	..	.	1001000....
4	00011	00	0	00.....	3.00	01001	..	.	00101100...
-7	00011	00	1	11.....	4.00	01001	..	.	10010000...
102.5	01011	01	0	0000000001..	7.00	01010	..	.	010111100..
1.025	01011	11	0	0000000001..	10.25	01011	..	.	0000000001..
総ビット	102 bits(13 bytes)				総ビット	93 bits(12 bytes)			

ろ、オリジナルデータは810,820の数値が、各10バイトで入力され、計8,108,200のバイトの容量となるが、この圧縮により得られたデータは1,065,453バイトであり13.14%という高圧縮を達成した。N式数値圧縮法は、表現桁数の大きい数値の比率が増えるほど圧縮率が下がり、

キャラクター形式の50%程度まで下がり、実質的な圧縮効果をもたなくなることも有り得るが、表現桁数の少ない数値の出現頻度もかなり存在する統計データの場合、おおよそ20%程度への圧縮は実現できる⁶⁾。

2. 国勢調査データファイルのデータベース化

1. で提案した「N式数値圧縮法」を利用し、国勢調査データのデータベース化を試みよう。同データファイルには、統計数値以外にも、各レコードの情報をしめすRID (Record Identifier) が入力されており、このレコード情報をどのようにインデックス化するかが、アクセス効率を確保するのに重要な問題となる。そこでこのRID情報がどのように入力されているか見てみよう (表-2 参照)。

同じ国勢調査データでも、提供されているデータファイルの種類により、RID部分はかなり異なっている。従って作成するデータベース・システムも、この相違を受けて個々にシステム構成することもできる。しかしそうすると、管理しなければならないファイルの数が多くなり、さらにデータベース検索ソフトが複雑になるし、検索手順もそれぞれについて異なり、ユーザーには操作しにくいものになる可能性が強い。インデックス・システムがある程度冗長性を持つことになっても、統一したデータ

ベース・システム構成にするのが、開発にも利用にも都合の良いことのように思える。

2-1. データベース・ファイルの構成

そこで以下のような3つのファイルからこのシステムを構成させる。データ検索者は通常、統計の種類、表番号という大きな分類項目から、都道府県、市区町村と順に小さく絞りながら情報を得ていくものという前提に立っており、重層的に構造化されたデータベース・システムとなっている。国勢調査基本集計を例にして内容を説明する。

2-1-1. メインファイル

統計数値データだけがレコードを単位としたN式数値圧縮法に従って入力されている(RCI-N1付)。1レコードのデータがバイトの途中で終わる時は、残りのビットには [0] b

表－3 国勢調査データファイルの RID 情報

基本集計		従業地・通学地		抽出詳細集計		調査区別集計		地域メッシュ統計	
セル数	5	セル数	5	セル数	5	セル数	5	ファイル名	6
ファイル名	5	ファイル名	5	ファイル名	5	ファイル名	5	ブランク	1
統計表番号	5	統計表番号	5	統計表番号	5	統計表番号	5	地域コード	9
都道府県	2	地域1区分	1	都道府県	2	ブランク	2	ブランク	2
市区町村	3	都道府県	2	市区町村	3	都道府県	2	区画区分	1
ブランク	2	市区町村	3	ブランク	2	市区町村	3	ブランク	1
人口集中地区	2	地域2区分	1	人口集中地区	2	ブランク	2	地形図名	18
分類項目1	3	都道府県	2	分類項目1	3	町丁・字	3	市区町村数	2
分類項目2	3	市区町村	3	分類項目2	3	ブランク	1	都道府県1	2
分類項目3	3	ブランク	3	分類項目3	3	調査区番号	7	市区町村1	3
ブランク	7	分類項目1	3	分類項目4	3	人口集中地区	2	府県市区2	5
		分類項目2	3	ブランク	4	都市計画区分	2	府県市区3	5
		ブランク	4			抽出層符号	6	府県市区4	5
						ブランク	2		
								府県市区8	5

数字は各項目が入力されているカラム数（バイト数）を示す

が付加されているので、レコードがバイトの途中から始まることはない。1レコードに含まれるデータ項目数（セル数）はこのファイルの中には与えられていないので、次の2つのファイルがないと、データを再現する（圧縮を解除する）ことはできない。なお、1都道府県データは統計の種類、表番号が同じなら基本的には連続して入力されており、その連続1単位がアドレスファイルの1レコードとなっている。

2－1－2．インデックスファイル

オリジナル・データファイルRIDの大部分の情報を保存する。市区町村コード [B04]，分類項目1～3（B05～B07），ダミー分類項目（B08～B11），該当するデータがメインファイルで記録されているバイト数 [REC]，の計9つの数値がN式整数（但，正負ビットと桁数制御ビット無し）の記録形式で入力されている。メインファイルのレコード数とこのファイルのレコード数は同じとなる。なお，ダミー分類項目には欠損値を入力している。

2－1－3．アドレスファイル

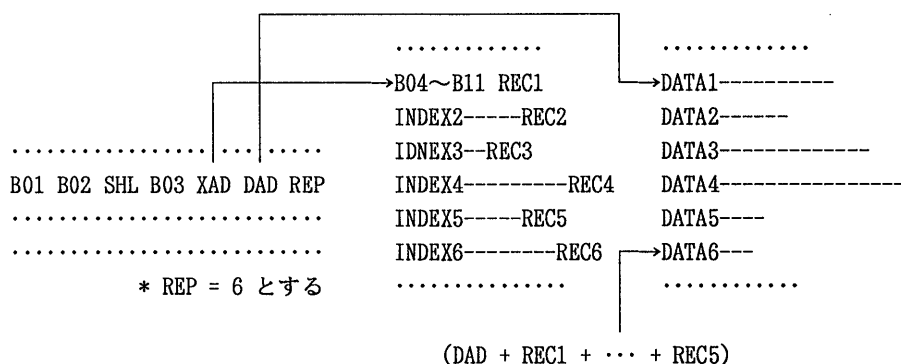
統計のファイル名を適当なコードに書換えた [B01]，表番号 [B02]，その表のデータ項目セル [SHL]，都道府県コード [B03]，該当するデータがインデックスファイルやメインファイルに入力されているレコードの総数 [REP]，インデックスファイルでのスタートアドレス [XAD]，データファイルのスタートアドレス [DAD] を入力している。このファイルは上記2ファイルの2％程度の容量となるので数値の圧縮は行わず，固定書式で入力されている。

2－1－4．データベースファイルの相互関係

この3つのファイル情報の関係を図示すると図－1のようになる。つまり，B01～B03の属性を持ったデータはSHL項目の数値データから構成されており，6ケースが入力されているが，DATA1は更にINDEX1に示されたB04～B11という属性を，DATA6はINDEX6に示される属性をあわせ持つ統計数値である。

このような構造になっているので，ファイル

図一1 データベース・ファイル情報の相互関係



内の情報を書き換えたり、特定のデータを削除しその領域をコンデンスすることなどは非常に難しいが、このシステムはデータベースの作成と利用に関わる2つの長所を持っている。それは情報を連結する（データベースを合併させる）ことが非常に簡単であることと、アドレスファイルをレコードを単位として自由に組み替えることができるということである⁷⁾。

磁気テープから資料—1で示すプログラムなどで変換作成された複数のデータベース・システムの結合は、年度や種類の差異にかかわらず、同じ操作で実行できる。インデックスとメインファイルに関しては、単純に連結コピーすればよいだけである（連結の順番は同じでなければならない）。アドレスファイルについては、連結コピーの際に前に置かれたインデックスとメインファイルの容量を、後ろに連結した側のアドレスファイルのXADとDADの値にそれぞれプラスすれば、あとはどのように連結コピーしてもよい。連結の際に何らかの情報操作をする必要があるのは、全体の2%程度の容量にしかならないアドレスファイルだけであり、磁気テープデータ、ファイルを入手の都度、データベース・システムを膨らますことができるという特性をもっている。

またアドレスファイルの各レコードは相互に独立したもののなので、レコード単位に並べ替えるのは全く自由であり、また、特定のレコードのみをとりだしたアドレスファイルを作成することも可能である。複数ユーザーに利用させる

環境で、ユーザー資格によって利用できる統計を限定したいときなどは、各ユーザーに対応したアドレスデータを抽出して利用させることも出来る。

2-2. データベースの作成プログラム

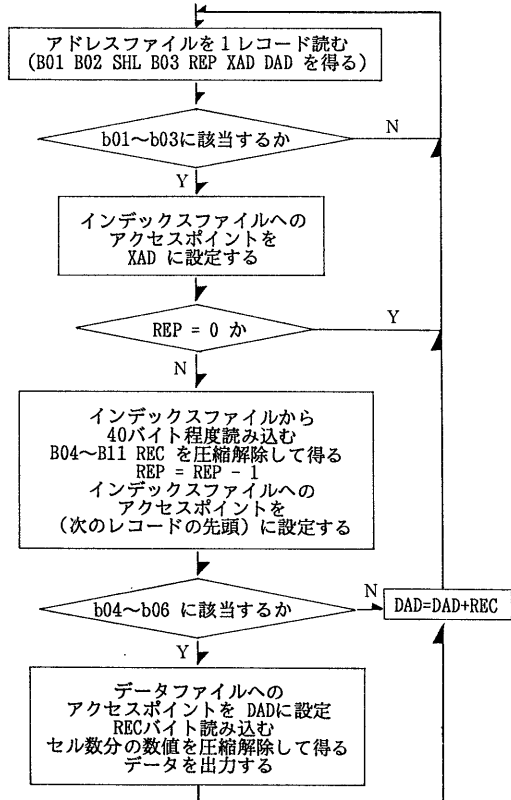
2-1. で示した3つのファイルから構成されるデータベース・ファイルを作成するPL/I言語によるプログラムは資料-1のようになる。これは国勢調査第1次～第3次基本集計のデータファイルを変換するものであるが、入力ファイル・変数宣言(11～49)と、一連の初期設定作業(54～63)をした後、入力データファイル1レコードを読み、(68)、分類の1～11(N_BUN(1)～N_BUN(11))、セル数(N_SHELL)、整数リスト(N_SUJI(1)～N_SUJI(N_SHELL))、その数値の繰り上がり桁数(N_KETA(1)～N_KETA(N_SHELL))に値を入れ(69～119)、サブプログラムNC_MAINをコールする(121)というものである。サブプログラムには、初期設定NC_INIT、メインルーチンNC_MAIN、最終処理NC_ENDの3つと、そこから更にコールされるもの(NC_COMP4, NC_UCOMP4, NC_STOCK, NC_MCOMP)があるが、すべてあわせても約230行とそれ程大きなものでないの、プログラム内に置いている。入力データファイルの書式と構成により変更する必要があるのは、資料-1で網かけされた部分(およそ70行)だけである。

2-3. データベースへのアクセス法

このデータベースへのアクセスは、検索したい統計種別 [b01], 表番号 [b02], 都道府県 [b03], 市区町村コード [b04], 分類1～分類3 [b05～b07] が指定されると、図-2のようになる。

このアクセスと検索を実行するC言語による(サブ)プログラムは資料-2に示したが、b01～b03は単一のコード指定、b04～b11はそれぞれについてコードリストを指定出来るようになっている。ユーザープログラムでの利用の仕方は、(1)最初にデータベースの3つのファイル名と欠損値に割当てた数値を指定してnc_open (84行)をコールする、(2)検索したい条件リストを指定してnc_set (103)をコールする、(3)出力結果を取込む変数のアドレス(ポインタ)を指定してncds (117)をコールするとい

図-2 データベース・アクセスのフローチャート



う3つの手順となる。同一条件に該当するデータをすべて取りだしたいなら、(3)を繰返し、別の条件でデータを取りだしたいなら(2)を再設定して(3)を繰返すという処理になる。

2-4. その他の統計資料のデータベース化

このシステムは、国勢調査データを対象にして設計したものであるが、結果として、書き換え型ではないたいの統計資料をこのデータベース・システムに組み込むことを可能とする構成になっている。数値データがある属性をもったレコードのような単位で取り扱うことができるもので、その属性を分類する項目を11項目以下にまとめることができるもの(さらに条件をつけるなら、大分類項目3種類以下で一定データ絞り込みが出来、中分類項目を8種類以下に整理できるもの)ならば、このデータベース・システムに組み込むことが可能である。特別に整理しなければならないのは分類項目1～11のコード対応表だけであり、後は、同一のソフト、同一の手順でデータを検索・出力させることができる。かなり広範囲に利用できる効率的な統計数値データベース・システムであると確信する。

付 記

このシステム構成に基づく国勢調査データベース・システム (BUDAS-POCE:Bukkyo University Data base System for POpulation CEnsus of japan) が完成している。出力形式は固定書式とLotus-1.2.3用K3フォーマットのみで、コメント文と都道府県・市区町村名以外は数字のみの出力である。統計表の種類は大学機関として購入したデータファイル量に制限されているが、京都府分(1975,80,85度分)が91年6月より学内の研究・教育用に利用できるようになる予定である。

このシステムの利用の中でユーザー・インターフェースの改良や、出力形式の改良(日本語ラベル付けなど)が求められるが、それらはプリ・ポストプロセッサソフトを適当に組み込

めばよいことであると理解しており、徐々に実行していくつもりである。

なお、このデータベース・システムの作成にあたって、立命館大学計算機センターには多くの便宜をはかって頂いた。深く感謝したい。

注

- 1) 排稿「大容量データベースの圧縮利用法——日経財務データの圧縮とパソコンでの利用」、『社会学部論叢』第23号、佛教大学学会、1989年12月。「家計調査年報のデータベース化」、同上、第24号、1990年12月。
- 2) ギルバート・ヘルド（渡部・堀池訳）、『データ圧縮技法入門』、啓学出版、1985年「訳者まえがき」より。
- 3) 例えば英語で書かれた文書の中に実際に表れる文字について考えてみると、出現頻度の高い文字とそうでないものがある。普通英字1文字は1バイト（8ビット）で表されるが、出現頻度の高いものは8ビットより短いコードで、低いものは場合によってはそれより多いビット例のコードで表現しても文書全体としてみれば圧縮出来るというものである。
- 4) パソコン通信の中で広まっている吉崎榮泰氏が作成したL H A R C法はデータの性質を問わない高圧縮を実現するものである。これは適応的ハフマン法を採用しており、圧縮の解除（解凍）はファイルをシーケンシャルにアクセスすることで行う。
- 5) コンピュータが扱う数値の記録形式は、2進数が都合の良いことは周知の事であるが、統計表（あるいは統計データファイル）は、それがどのように計算された数値であるにせよ、10進数で

表記・記録されるのが一般的である。この数値が整数の場合は問題ないが、小数点以下表示がある場合は直接2進数で表記するとかなり長大なビット列が必要となることがあり、10進数の小数点以下を繰上げた方が、簡略に表現出来ることが多い。

- 6) N式数値圧縮法はさらに桁数インデックスに使用するビット数とそのコード化について「改良」することができる。すなわち表-1(1)を見ても各統計書によって、表現必要ビット数の分布は欠損値を別とすると、8ビット、もしくは10ビットを頂点とするような分布となっているので、表現必要ビット数を現すコードをハフマンの記号化法に従って $-\text{LOG}_2(P)$ （Pはその表現必要ビット数の数値出現頻度）ビットに設定することも出来る。また出現頻度が1%を越えるような数値（1～9など）を、この記号化に直接組み込み、データ表現のためのビットを省略することもできる。このような「改良」の可能性を認識しながら、それを採用しないのは圧縮のためのコード化を、データによって変動させることが、データベースの追記性の障害となることを避けようとしたからである。データベース・システムの運用効率を高めるためのある程度の冗長性は避けられないものと考えたからである。
- 7) このような「長所」は結果的なものでなく、最初の設計で意図したことでもある。この調査のデータファイルは、集計が完了してから順に公表されるものであり、また我々がデータを入手するのは、必要に応じた地域の統計からである。どのような集め方をしても、結果として同じ手順で利用できるデータ・ベースシステムにする必要があったわけである。

資料1 国勢調査データベース・システム作成プログラム

```

1 BUDASP : PROC OPTIONS(MAIN) ;
2 /* ***** MAKE BUDAS-POCE ***** */
3 *
4 * 国勢調査・基本集計データファイルの圧縮
5 *
6 * 1990-10-24 (REV-91-02-20) 内藤三義
7 *
8 *****
9
10 ファイル宣言
11 DCL INF FILE RECORD INPUT ;
12 DCL OUF FILE RECORD OUTPUT ;
13 DCL IDX FILE RECORD OUTPUT ;
14 DCL ADR FILE RECORD OUTPUT ;
15 DCL ENV(VB RECSIZE(1000) BLKSIZE(4000)) ;
16 DCL ENV(VB RECSIZE(48) BLKSIZE(2400)) ;
17 DCL ENV(FB RECSIZE(24) BLKSIZE(2400)) ;
18
19 /* =====
20 変数宣言
21 DCL 1 IDATA,
22 2 ID CHAR(32000) VARYING ;
23 DCL 1 IDATA,
24 2 SHELL FIXED BIN(15), 2 NUM FIXED BIN(30),
25 2 F FIXED BIN(8), 2 TENKEN CHAR(5),
26 2 SUJ11 FIXED BIN(15), 2 SUJ12 FIXED BIN(15),
27 2 C10 CHAR(10), 2 C5 CHAR(5),
28 2 C3 CHAR(3), 2 C2 CHAR(2) ;
29 DCL 1 NCDATA,
30 2 NCID,
31 3 N_BUN(12) FLOAT BIN(53), 3 N_SHELL FLOAT BIN(53),
32 3 N_VAR(2000) FLOAT BIN(53),
33 3 N_KET(2000) FIXED BIN(8),
34 2 N_MIS FLOAT BIN(53) INIT(-999999999),
35 2 N_OD CHAR(2048) VARYING,
36 2 NCOA,
37 3 N_OABUN1 CHAR(2), 3 N_OABUN2 CHAR(4),
38 3 N_OASHEL CHAR(2), 3 N_OABUN3 CHAR(4),
39 3 N_OAXADR CHAR(4), 3 N_OAXREP CHAR(4),
40 3 N_OAADR CHAR(4),
41 2 NCPARM,
42 3 N_ORGB(4) FLOAT BIN(53), 3 N_NUM FIXED BIN(30),
43 3 N_KER BIN(4), 3 N_KEP FIXED BIN(8),
44 3 N_PMP BIN(1), 3 N_REC FIXED BIN(15),
45 3 N_XADR FLOAT BIN(53), 3 N_XREP FIXED BIN(31),
46 3 N_DADR FLOAT BIN(53), 3 N_XMAX FIXED BIN(15),
47 3 N_DMAX FIXED BIN(15),
48 3 N_C4 CHAR(4), 3 N_BP FIXED BIN(15),
49 3 N_SUJ1 FLOAT BIN(53), 3 N_BS BIT(64) ;
50
51 /* =====
52 初期設定とファイル・オープン
53 CALL NC_INIT ;
54 OPEN FILE (INF), FILE (OUF), FILE (IDX), FILE (ADR) ;
55 ON ENDFILE (INF) GOTO FINISH ;
56 PUT EDIT ('INPUT FILE NAME OF CENSUS ( 5 CHARA ) --> ')
57 (SKIP,A(42)) ;
58 GET EDIT (TENKEN) (A(5)) ;
59 PUT EDIT ('INPUT SENSUS NO. ( UNDER 65535 ) -----> ')
60 (SKIP,A(42)) ;
61 GET LIST (SUJ11) ;
62 N_BUN(1) = SUJ11 ;
63 N_BUN(9) = N_MIS; N_BUN(10) = N_MIS; N_BUN(11) = N_MIS ;
64
65 /* =====
66 読み込み開始
67
68 START : READ FILE (INF) INTO (ID) ;
69 NUM = NUM + 1 ;
70 OSIZE = OSIZE + LENGTH(ID) ;
71 IF SUBSTR(ID,6,5) = TENKEN THEN GOTO STEP1 ;
72 PUT EDIT ('DATA ERROR CASE NO.--> ',NUM)
73 (SKIP,A(25),F(7)) ;
74 GOTO START ;
75
76 /* アドレスデータ (分類 2, 3, シェル) の取込み
77 STEP1 :
78 C5 = SUBSTR(ID,1,5) ; IF C5 = ' ' THEN C5 = ' 0' ;
79 N_SHELL = C5 ; SHELL = C5 ;
80 C3 = SUBSTR(ID,11,3) ; IF C3 = ' ' THEN C3 = ' 0' ;
81 SUJ11 = C3 ; SUJ11 = SUJ11 * 100 ;
82 C2 = SUBSTR(ID,14,2) ; IF C2 = ' ' THEN C2 = ' 0' ;
83 SUJ12 = C2 ; N_BUN(2) = SUJ11 + SUJ12 ;
84 C2 = SUBSTR(ID,16,2) ; IF C2 = ' ' THEN C2 = ' 0' ;
85 N_BUN(3) = C2 ;
86
87 /* インデックス (分類 4~8) の取込み
88 STEP2 :
89 C3 = SUBSTR(ID,18,3) ; IF C3 = ' ' THEN C3 = ' 0' ;
90 N_BUN(4) = C3 ;
91 C2 = SUBSTR(ID,23,2) ;
92 IF C2 = '00' THEN N_BUN(5) = 1 ;
93 ELSE N_BUN(5) = 0 ;
94 DO I = 1 TO 3 ;
95 C3 = SUBSTR(ID,I*3+22,3) ;
96 IF C3 = ' ' THEN N_BUN(I+5) = N_MIS ;
97 ELSE N_BUN(I+5) = C3 ;
98 END ;
99
100 /* 実データの取込み
101 STEP3 :
102 DO I = 1 TO SHELL ;
103 C10 = SUBSTR(ID,I*10+31,10) ;
104 IF C10 = ' ' THEN GOTO STE31 ;
105 IF C10 = ' ' THEN GOTO STE31 ;
106 IF C10 = ' ' THEN GOTO STE31 ;
107 DO WHILE (SUBSTR(C10,10,1) = ' ') ;
108 C10 = ' ' !! SUBSTR(C10,1,9) ;
109 END ;
110 P = INDEX(C10,'.') ;
111 IF P = 0 THEN N_KET(I) = 0 ;
112 ELSE DO ;
113 N_KET(I) = 10 - P ;
114 SUBSTR(C10,1,P) = ' ' !! SUBSTR(C10,1,P-1) ;
115 END ;
116 N_VAR(I) = C10 ;
117 GOTO STE32 ;
118 STE31 : N_VAR(I) = N_MIS ; N_KET(I) = 0 ;
119 STE32 : END ;
120
121 CALL NC_MAIN ;
122
123 GOTO START ;
124
125 /* =====
126 SUB ROUTINE
127
128 /* 初期化
129 NC_INIT : PROC ;
130 N_XADR = 0 ; N_DADR = 0 ; N_XREP = 0 ;
131 N_NUM = 0 ; N_DMAX = 0 ; N_XMAX = 0 ; N_BUN = N_MIS ;
132 END NC_INIT ;
133
134 /* メイン・サブルーチン
135 NC_MAIN : PROC ;
136 DCL ( NC_P, NC_M ) FIXED BIN(15),
137 NC_BT6 BIT(6), NC_BT4 BIT(4) ;
138
139 /* ..... 異なる表などの時アドレスの最終計算と書きだし .....
140 N_NUM = N_NUM + 1 ;
141 IF N_NUM = 1 THEN GOTO NC_STEP1 ;
142
143 IF N_BUN(1) = N_ORGB(1) & N_BUN(2) = N_ORGB(2)
144 & N_BUN(3) = N_ORGB(3) THEN GOTO NC_STEP2 ;
145 N_SUJ1 = N_XREP ; CALL NC_UCOMP4 ; N_OAXREP = N_C4 ;
146 N_XREP = 0 ;
147 WRITE FILE (ADR) FROM(NCOA) ;
148
149 /* ..... 異なる表などの時当たらしアドレスの計算 .....
150 NC_STEP1 : N_ORGB(1) = N_BUN(1) ; N_ORGB(2) = N_BUN(2) ;
151 N_ORGB(3) = N_BUN(3) ; N_ORGB(4) = N_SHELL ;
152 N_SUJ1 = N_BUN(1) ; CALL NC_UCOMP4 ;
153 N_OABUN1 = SUBSTR(N_C4,1,2) ;
154 N_SUJ1 = N_BUN(2) ; CALL NC_UCOMP4 ; N_OABUN2 = N_C4 ;
155 N_SUJ1 = N_SHELL ; CALL NC_UCOMP4 ;
156 N_OASHEL = SUBSTR(N_C4,1,2) ;
157 N_SUJ1 = N_BUN(3) ; CALL NC_UCOMP4 ; N_OABUN3 = N_C4 ;
158 N_SUJ1 = N_XADR ; CALL NC_UCOMP4 ; N_OAXADR = N_C4 ;
159 N_SUJ1 = N_DADR ; CALL NC_UCOMP4 ; N_OAADR = N_C4 ;
160
161 /* ..... 同一表でセル数が異なれば WARNING .....
162 NC_STEP2 : N_XREP = N_XREP + 1 ;
163 IF N_SHELL = N_ORGB(4) THEN PUT EDIT
164 ('*WARNING* SHELL DATA ERROR CASE NO. TRUE -> REAL',
165 '[',N_NUM,'] [' ,N_ORGB(4),
166 ']->[' ,N_SHELL,']')
167 (SKIP,A(51),SKIP,A(27),F(7),A(3),F(5),A(4),F(5),A(1)) ;
168
169 /* ..... 桁数インデックスのためのパラミータの取得 .....
170 NC_STEP3 : N_KEP = -1 ; N_KET = -1 ; N_PMP = 0 ;
171 NC_P = 0 ; NC_M = 0 ;
172 DO I = 1 TO N_SHELL ;
173 IF N_VAR(I) = N_MIS THEN DO ;
174 IF N_KER = N_KET(1) THEN N_KEP = N_KEP + 1 ;
175 IF N_KER < N_KET(1) THEN N_KER = N_KET(1) ;
176 IF N_VAR(I) > 0 THEN NC_P = NC_P + 1 ;
177 IF N_VAR(I) < 0 THEN NC_M = NC_M + 1 ;
178 END ;
179 END ;
180
181 IF N_KEP = -1 THEN N_KEP = 0 ;
182 IF N_KER = -1 THEN N_KER = 0 ;
183
184 /* ..... 桁数インデックス RCI-NO の決定 .....
185 NC_STEP4 : N_REC = 0 ; N_OD = 0 ; N_BP = 1 ;
186 SUBSTR(N_BS,1,3) = '111' ; N_BP = 4 ;
187 IF NC_M = 0 THEN DO ;
188 SUBSTR(N_BS,N_BP,1) = '0' ; N_BP=N_BP+1 ; N_PMP=0 ;

```

```

189 END ;
190 ELSE IF NC_P = 0 THEN DO ;
191   SUBSTR(N_BS, N_BP, 2) = '10' ; N_BP = N_BP + 2 ; N_PMP = 0 ;
192 END ;
193 ELSE DO ;
194   SUBSTR(N_BS, N_BP, 2) = '11' ; N_BP = N_BP + 2 ; N_PMP = 1 ;
195 END ;
196 NC_BT4 = N_KER ;
197 IF N_KEP THEN DO ;
198   SUBSTR(N_BS, N_BP, 1) = '1' ;
199   N_KER = 5 - INDEX(NC_BT4, '1') ; NC_BT4 = N_KER ;
200 END ;
201 ELSE SUBSTR(N_BS, N_BP, 1) = '0' ;
202 N_BP = N_BP + 1 ;
203 SUBSTR(N_BS, N_BP, 4) = NC_BT4 ; N_BP = N_BP + 4 ;
204
205 /* .... 桁数インデックスを RCI-N1 への変形 ..... */
206 NC_STEP5 :
207   NC_BT6 = SUBSTR(N_BS, 4, 6) ;
208   IF NC_BT6 = '000000' THEN DO ;
209     SUBSTR(N_BS, 1, 2) = '00' ; N_BP = 3 ;
210   END ;
211   IF NC_BT6 = '010010' THEN DO ;
212     SUBSTR(N_BS, 1, 2) = '01' ; N_BP = 3 ;
213   END ;
214   IF NC_BT6 = '000010' THEN DO ;
215     SUBSTR(N_BS, 1, 3) = '100' ; N_BP = 4 ;
216   END ;
217   IF NC_BT6 = '010001' THEN DO ;
218     SUBSTR(N_BS, 1, 3) = '101' ; N_BP = 4 ;
219   END ;
220   IF NC_BT6 = '000001' THEN DO ;
221     SUBSTR(N_BS, 1, 3) = '110' ; N_BP = 4 ;
222   END ;
223
224 /* .... 実データ部分の圧縮書きだし ..... */
225 NC_STEP6 :
226   DO I = 1 TO N_SHELL ;
227     N_KER = N_KEP(I) ;
228     N_SUJI = N_VAR(I) ; CALL NC_MCOMP ; CALL NC_STOCK ;
229   END ;
230   IF N_BP > 1 THEN DO ;
231     SUBSTR(N_BS, N_BP, 7) = '00000000' ; N_BP = 9 ;
232     CALL NC_STOCK ;
233   END ;
234   IF N_DMAX < N_REC THEN N_DMAX = N_REC ;
235   N_DADR = N_DADR + N_REC ;
236   WRITE FILE (OUF) FROM (N_OD) ;
237
238 /* .... インデックスデータの最終計算と書きだし ..... */
239 NC_STEP7 : N_KEP = 0 ; N_PMP = 0 ; N_OD = '' ;
240 N_BP = 1 ; N_BUN(12) = N_REC ; N_REC = 0 ;
241 DO I = 4 TO 12 ;
242   N_SUJI = N_BUN(I) ; CALL NC_MCOMP ; CALL NC_STOCK ;
243 END ;
244 IF N_BP > 1 THEN DO ;
245   SUBSTR(N_BS, N_BP, 7) = '00000000' ; N_BP = 9 ;
246   CALL NC_STOCK ;
247 END ;
248 IF N_XMAX < N_REC THEN N_XMAX = N_REC ;
249 N_XADR = N_XADR + N_REC ;
250 WRITE FILE (IDX) FROM (N_OD) ;
251 END NC_MAIN ;
252
253 /* 最終アドレスデータの書きだし ..... */
254 NC_END : PROC ;
255   N_SUJI = N_XREP ; CALL NC_UCOMP4 ; N_OAXREP = N_C4 ;
256   WRITE FILE (ADR) FROM (NCOA) ;
257   N_SUJI = N_XADR ; CALL NC_UCOMP4 ; N_OAXADR = N_C4 ;
258   N_SUJI = N_DADR ; CALL NC_UCOMP4 ; N_OADADR = N_C4 ;
259   UNSPEC(N_OABUN1) = 'FFFFFF' B4 ;
260   UNSPEC(N_OABUN2) = 'FFFFFF' B4 ;
261   UNSPEC(N_OASHL) = 'FFFFFF' B4 ;
262   UNSPEC(N_OABUN3) = 'FFFFFF' B4 ;
263   UNSPEC(N_OAXREP) = 'FFFFFF' B4 ;
264   WRITE FILE (ADR) FROM (NCOA) ;
265
266   PUT EDIT (' TOTAL RECORD NUMBER=', N_NUM)
267   (SKIP, A(26), F(15)) ;
268   PUT EDIT (' FILE SIZE OF DATA=', N_DADR)
269   (SKIP, A(26), F(15)) ;
270   PUT EDIT (' INDEX=', N_XADR)
271   (SKIP, A(26), F(15)) ;
272   PUT EDIT (' MAX RECORD SIZE OF DATA=', N_DMAX)
273   (SKIP, A(26), F(15)) ;
274   PUT EDIT (' INDEX=', N_XMAX)
275   (SKIP, A(26), F(15)) ;
276 END NC_END ;
277
278 /* -----
279 SUB-SUB ROUTINE
280 ----- */
281

```

```

282 /* .... GET LONG (4 BYTES) .. THIS PROCESS IS NOT USED */
283 NC_COMP4 : PROC ;
284   DCL NC_BT32 BIT(32), NC_BN31 FIXED BIN(31) ;
285   NC_BN31 = N_SUJI ; NC_BT32 = UNSPEC(NC_BN31) ;
286   UNSPEC(HC4) = SUBSTR(NC_BT32, 25, 8) !! SUBSTR(NC_BT32, 17, 8)
287   !! SUBSTR(NC_BT32, 9, 8) !! SUBSTR(NC_BT32, 1, 8) ;
288 END NC_COMP4 ;
289
290 /* .... GET UNSIGNED LONG (4 BYTES) ..... */
291 NC_UCOMP4 : PROC ;
292   DCL NC_BT32 BIT(32), NC_BT64 BIT(64),
293   NC_BN7 FIXED BIN(7) ;
294   NC_BT64 = UNSPEC(N_SUJI) ;
295   NC_BN7 = SUBSTR(NC_BT64, 2, 7) ;
296   NC_BT32 = '00000000' B4 ;
297   IF NC_BN7 = 0 THEN UNSPEC(N_C4) = NC_BT32 ;
298   ELSE DO ;
299     NC_BN7 = (NC_BN7 - 64) * 4 ;
300     SUBSTR(NC_BT32, 33 - NC_BN7, NC_BN7) =
301     SUBSTR(NC_BT64, 9, NC_BN7) ;
302     UNSPEC(N_C4) = SUBSTR(NC_BT32, 25, 8)
303     !! SUBSTR(NC_BT32, 17, 8) !! SUBSTR(NC_BT32, 9, 8)
304     !! SUBSTR(NC_BT32, 1, 8) ;
305   END ;
306 END NC_UCOMP4 ;
307
308 /* .... STOCK ONE CHARACTER IN [N_OD] ..... */
309 NC_STOCK : PROC ;
310   DCL NC_C1 CHAR(1), NC_BT8 BIT(8) ;
311   DO WHILE (N_BP > 8) ;
312     NC_BT8 = SUBSTR(N_BS, 1, 8) ; UNSPEC(NC_C1) = NC_BT8 ;
313     N_BS = SUBSTR(N_BS, 9, 56) ; N_BP = N_BP - 8 ;
314     N_OD = N_OD !! NC_C1 ; N_REC = N_REC + 1 ;
315   END ;
316 END NC_STOCK ;
317
318 /* .... MAKE ONE COMPRESS DATA ON N_BS ..... */
319 NC_MCOMP : PROC ;
320   DCL NC_BT64 BIT(64), NC_BT56 BIT(56), NC_BT7 BIT(7),
321   NC_BT4 BIT(4), NC_BT1 BIT(1), NC_P BIN(8),
322   NC_BN7 FIXED BIN(7) ;
323   IF N_SUJI = N_MIS THEN GOTO NC_MSTEL1 ;
324   SUBSTR(N_BS, N_BP, 2) = '11' ; N_BP = N_BP + 2 ;
325   GOTO NC_MSTEL2 ;
326   NC_MSTEL1 : NC_BT64 = UNSPEC(N_SUJI) ;
327   NC_BT1 = SUBSTR(NC_BT64, 1, 1) ;
328   NC_BN7 = SUBSTR(NC_BT64, 2, 7) ;
329   IF NC_BN7 = 0 THEN DO ;
330     NC_BN7 = (NC_BN7 - 64) * 4 ;
331     NC_BT56 = SUBSTR(NC_BT64, 9, 56) ;
332     NC_P = INDEX(NC_BT56, '1') ;
333     NC_BN7 = NC_BN7 - NC_P + 1 ;
334   END ;
335   NC_BT7 = NC_BN7 ;
336   IF NC_BN7 < 16 THEN DO ;
337     SUBSTR(N_BS, N_BP, 5) = SUBSTR(NC_BT7, 3, 5) ;
338     N_BP = N_BP + 5 ;
339   END ;
340   ELSE DO ;
341     SUBSTR(N_BS, N_BP, 2) = '10' ; N_BP = N_BP + 2 ;
342     SUBSTR(N_BS, N_BP, 4) = SUBSTR(NC_BT7, 4, 4) ;
343     N_BP = N_BP + 4 ;
344   END ;
345   IF N_KEP THEN DO ;
346     NC_BT4 = N_KER ;
347     SUBSTR(N_BS, N_BP, N_KEP) =
348     SUBSTR(NC_BT4, 5 - N_KEP, N_KEP) ;
349     N_BP = N_BP + N_KEP ;
350   END ;
351   IF NC_BN7 = 0 THEN GOTO NC_MSTEL2 ;
352   IF N_PMP THEN DO ;
353     SUBSTR(N_BS, N_BP, 1) = NC_BT1 ; N_BP = N_BP + 1 ;
354   END ;
355   IF NC_BN7 = 1 THEN GOTO NC_MSTEL2 ;
356   NC_BN7 = NC_BN7 - 1 ;
357   SUBSTR(N_BS, N_BP, NC_BN7) =
358   SUBSTR(NC_BT56, NC_P + 1, NC_BN7) ;
359   N_BP = N_BP + NC_BN7 ;
360   NC_MSTEL2 : END NC_MCOMP ;
361
362 /* ===== END SUBROUTINE ===== */
363
364 FINISH :
365   PUT EDIT ('END PROCESS ORIGINAL SIZE=', OSIZE)
366   (SKIP, A(26), F(15)) ;
367   CALL NC_END ;
368   CLOSE FILE (INF), FILE (OUF), FILE (IDX), FILE (ADR) ;
369   END BUDASP ;

```

資料2 N式社会統計データベース検索出力プログラム

```

1  /* =====
2
3  Naito's Social Statistics Compressed Database System
4
5  N式社会統計圧縮データベース・システム
6
7  ncds.c v1.0 1991-01-15 内藤 三義
8  ===== */
9
10 #include <stdio.h>
11 #include <conio.h> /* kbhit getch */
12 #include <process.h> /* exit */
13
14 /*
15  nc_open(fln,mv) NCDS ファイルオープンと欠損値の指定
16  char *(*fln) double mv
17  fin : アドレスファイル・パス名
18  fin+1 : インデックスファイル・パス名
19  fin+2 : データファイル・パス名
20  mv : 欠損値に割り当てる数値
21  return = 0 : 正常終了
22  = 1 : アドレスファイルがオープン出来ない
23  = 2 : インデックスファイルがオープン出来ない
24  = 4 : データファイルがオープン出来ない
25
26  nc_close() NCDS ファイルのクローズ
27  return値なし
28
29  nc_set(list,param) NCDS 検索条件の設定
30  long *(list) int param
31  (1) 検索したい分類1~11のコード(リスト)の
32  先頭アドレスを list~list+11 に設定する
33  (2) param = 0 : 現在のアドレスポイントから検索
34  = 1 : アドレスファイルの最初から検索する
35
36  return = 0 : 正常終了
37  = 1 : ファイルがオープンされてない
38
39  ncds(data,keta,param) NCDS 検索と出力
40  double *data char *keta int param
41  (1) data : 出力する数値を収めるポインタ
42  data[0] : レコード番号
43  data[1] : セル数
44  data[2~4] : 分類1~3
45  data[5~7] : xadr, xrep, dadr
46  data[8] : 相対レコード番号
47  data[9~16] : 分類4~11
48  data[17~*] : データ・リスト
49  data[*+1] : 欠損値
50  (2) keta : 出力する数値を収めるポインタ
51  keta[0~*] : 対応する繰上り桁数
52  keta[*+1] : -1
53  (3) param = 1 : アドレス情報のみ出力
54  = 2 : アドレス・インデックス情報を出力
55  = 3 : アドレス・インデックス・データ情報
56  を出力
57
58  return = 0 : 正常終了
59  = -1 : アドレスファイル・エラー
60  = -2 : インデックスファイル・エラー
61  = -3 : データファイル・エラー
62  = 1 : 対応データはない(eof of address file)
63  = 2 : 検索条件が指定されていない
64  = 3 : param に1~3 以外が指定された
65
66  struct { unsigned int bun1; unsigned long bun2;
67  unsigned int shell;
68  unsigned long bun3, xadr, repeat, dadr;
69  } s_adr;
70 FILE *fp_a, *fp_x, *fp_d;
71 double d_missing, d_xdat[12], *dp_1, *dp_2, *dp_3;
72 long l_anum, l_xnum, l_adat[9];
73 int i_fflag = 0, i_xdp, i_xdl, i_dparm = -1,
74 i_lparm = -1;
75 char c_xd[2200], c_dd[2048], *cp_1;
76
77 /* ===== */
78 int nc_open(char *(*fln), double mv)
79 { int flag = 0;
80 if ((fp_a = fopen(*fln, "rb")) != NULL) flag++;
81 if ((fp_x = fopen(*fln, "rb")) != NULL) flag += 2;
82 if ((fp_d = fopen(*fln, "rb")) != NULL) flag += 4;
83 d_missing = mv; i_fflag = flag;
84 flag = 7 - flag;
85 if (flag) nc_close();
86 else { l_anum = 0; l_xnum = -1; }
87 return(flag);
88 }
89
90 /* ===== */
91 nc_close()
92 {
93 }
94
95 96 { if ((i_fflag & 1) == 1) fclose(fp_a);
97 if ((i_fflag & 2) == 2) fclose(fp_x);
98 if ((i_fflag & 4) == 4) fclose(fp_d);
99 i_fflag = 0;
100 }
101 /* ===== */
102 int nc_set(long *(list), int parm)
103 { int i, flag = 0;
104 for (i = 0; i < 11; i++)
105 { lp_list[i] = *list; list++; }
106 i_lparm = 0;
107 if (parm == 1)
108 { if (i_fflag == 7)
109 { rewind(fp_a); l_anum = 0; l_xnum = -1;
110 i_dparm = -1; }
111 else { flag = 1; return(flag); }
112 }
113 return(flag);
114 }
115 /* ===== */
116 int ncds(double *data, char *keta, int parm)
117 { int i, j, flag;
118 if (i_fflag != 7)
119 { flag = 7 - i_fflag; return(flag); }
120 if (i_lparm)
121 { flag = 2; return(flag); }
122 dp_1 = data; cp_1 = keta;
123 if (parm == 1)
124 { flag = adrcheck(); return(flag); }
125 dp_2 = dp_1; dp_2 += 8;
126 if (parm == 2)
127 { flag = idxcheck(); return(flag); }
128 dp_3 = dp_1; dp_3 += 17;
129 if (parm == 3)
130 { int rec, shel, fsp = 0, op = 0;
131 if (i_dparm)
132 { if (flag == idxcheck()) return(flag); }
133 rec = (int)d_xdat[9]; shel = (int)l_adat[1];
134 flag = -3; i_dparm = -1;
135 if (fseek(fp_d, l_adat[7], fsp) != 0)
136 return(flag);
137 if (fread(&c_dd[0], rec, 1, fp_d) != 1)
138 return(flag);
139 rcomp(&c_dd[0], dp_3, cp_1, shel, op);
140 flag = 0; return(flag);
141 }
142 flag = 3; return(flag);
143 }
144
145 /* ----- アドレスファイルのチェック ----- */
146 return -1:file error, 1:end of file
147 }
148 /* ===== */
149 int adrcheck()
150 { int i, flag, rec = 24;
151 double *dp;
152 dp = dp_1; l_xnum = -1; i_dparm = -1;
153 while((fread((char *)&s_adr, rec, 1, fp_a) == 1))
154 { l_anum++;
155 if (s_adr.bun1 == 0xffff)
156 { flag = 1; return(flag); }
157 l_adat[2] = (long)s_adr.bun1; l_adat[3] = s_adr.bun2;
158 l_adat[4] = (long)s_adr.bun3;
159 flag = 0;
160 for (i = 0; i < 3; i++)
161 { if ((lp_list[i] != -1) &&
162 (*lp_list[i] != l_adat[2 + i])) flag++;
163 }
164 if (flag == 0)
165 { l_adat[0] = l_anum;
166 l_adat[1] = (long)s_adr.shell;
167 l_adat[2] = (long)s_adr.bun1;
168 l_adat[3] = s_adr.bun2;
169 l_adat[4] = (long)s_adr.bun3;
170 l_adat[5] = s_adr.xadr;
171 l_adat[6] = s_adr.repeat;
172 l_adat[7] = s_adr.dadr;
173 l_xnum = 0;
174 for (i = 0; i < 8; i++)
175 { *dp = (double)l_adat[i]; dp++; }
176 *dp = -1; return(flag);
177 }
178 flag = -1; return(flag);
179 }
180 }
181 /* ----- インデックス・ファイルのチェック ----- */
182 return -2:file error, 2:end of one address cluster
183 }
184 int idxcheck()

```

```

180 { int i, j, flag, fsp = 0, rec = 2048, cx = 9, op = 1;
181 long temp;
182 double *dp;
183 flag = 1; dp = dp_2; i_dparm = -1;
184 while(flag)
185 { if (l_xnum == -1)
186 { if (flag == adrcheck()) return( flag ); }
187 if (l_xnum == 0)
188 { if (fseek(fp_x, l_adat[5], fsp) != 0)
189 { flag = -2; return( flag ); }
190 fread(&c_xd[0], rec, 1, fp_x);
191 i_xdp = 0; i_xdl = 2048;
192 }
193 }
194 else
195 { l_adat[7] += (long)d_xdat[8];
196 if (i_xdp > (i_xdl - 48))
197 { j = 0;
198 for (i = i_xdp; i < 2048; i++)
199 { c_xd[j] = c_xd[i]; j++; }
200 xdp = 0; i_xdl = 2048 + j;
201 fread(&c_xd[j], rec, 1, fp_x);
202 }
203 }
204 }
205 l_xnum++; d_xdat[0] = (double)l_xnum;
206 i_xdp +=
207 rcomp(&c_xd[i_xdp], &d_xdat[1], cp_1, cx, op );
208 for (i = 1; i < 9; i++)
209 { lp_1 = lp_list[i + 2]; flag = 0;
210 if (*lp_1 != -1)
211 { flag = 1; temp = (long)d_xdat[i];
212 while((*lp_1 != -1) && (flag))
213 { if (temp == *lp_1) flag = 0; lp_1++; }
214 }
215 if (flag) i = 9;
216 if (l_xnum == l_adat[6]) l_xnum = -1;
217 }
218 for (i = 0; i < 9; i++) { *dp = d_xdat[i]; dp++; }
219 *dp = -1; i_dparm = 0;
220 return( flag );
221 }
222
223 /* ..... N式数値圧縮法からの復帰 ..... */
224 int rcomp( char *si, double *dl, char *bx, int cx, int i )
225 { char ah, al, bl, bh, cl, ch, dl, dh;
226 double ax;
227 int rec = 0;
228
229 /* .....
230 . al : one character cl : dummy
231 . ah : keta control bit ch : keta su
232 . bl : current attack bit(0-7) dl : databit length
233 . bh : +- control parameter dh : +- sign
234 . di : address of output data(long)
235 . si : address of input data(char)
236 . bx : address of output data for keta(char)
237 ..... */
238
239 dh = ah = bh = cl = 0;
240
241 if (i == 1) { ah = bh = dh = dl = ch = 0; bl = 7; }
242 else
243 { al = *si; si++; rec++; bl = 1;
244 if (al < 0) cl = 2;
245 al *= 2;
246 if (al < 0) cl++;
247 if (cl == 0) { ah = bh = dh = ch = 0; }
248 else if (cl == 1) { ah = 2; bh = dh = 0; }
249 else
250 { cl *= 2; bl++; al *= 2;
251 if (al < 0) cl++;
252 if (cl == 4) { ah = bh = dh = 0; ch = 2; }
253 else if (cl == 5) { ah = 1; bh = dh = 0; }
254 else if (cl == 6) { ah = bh = dh = 0; ch = 1; }
255 else cl = -1;
256 }
257 }
258 if (cl < 0)
259 { cl = 0; al *= 2; bl++;
260 if (al < 0)
261 { al *= 2; bl++;
262 if (al < 0) bh = 1;
263 else { bh = 0; dh = 1; }
264 }
265 }
266
267
268
269
270
271
272
273
274

```

```

275 else { bh = 0; dh = 0; }
276 al *= 2; bl++;
277 if (al < 0) ah = 1;
278 else ah = 0;
279 cl = 0;
280 for (i = 0; i < 4; i++)
281 { bl++; al *= 2; cl *= 2;
282 if (bl == 8) { al = *si; si++; rec++; bl = 0; }
283 if (al < 0) cl++;
284 }
285 if (ah) ah = cl;
286 else ch = cl;
287
288 if (ah == 0)
289 for (i = 0; i < cx; i++) { *bx = ch; bx++; }
290
291 /* .....
292 . bh = 1 : 正負制御ビットあり
293 . bh = 0 なら ( dh = 0 : 全て正 / dh = 1 : 全て負 )
294 . ah = 桁数制御ビット数 ( = 0 はなし )
295 ..... */
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357

```

実データの圧縮解除

```

358 while (cx > 0)
359 { cx--;
360 dl = 0;
361 for (i = 0; i < 2; i++)
362 { bl++; al *= 2; dl *= 2;
363 if (bl == 8)
364 { al = *si; si++; rec++; bl = 0; }
365 if (al < 0) dl++;
366 }
367 if (dl == 3)
368 { ax = d_missing; if (ah) { *bx = 0; bx++; } }
369 else
370 {
371 /* ..... 実データ記録 bit 数の取得 ..... */
372 cl = 3; if (dl == 2) cl = 4;
373 for (i = 0; i < cl; i++)
374 { bl++; al *= 2; dl *= 2;
375 if (bl == 8)
376 { al = *si; si++; rec++; bl = 0; }
377 if (al < 0) dl++;
378 }
379 if (dl > 16) dl -= 16;
380
381 /* ..... 小数点以下桁数の取得 ..... */
382 if (ah)
383 { ch = 0;
384 for (i = 0; i < ah; i++)
385 { bl++; al *= 2; ch *= 2;
386 if (bl == 8)
387 { al = *si; si++; rec++; bl = 0; }
388 if (al < 0) ch++;
389 }
390 *bx = ch; bx++;
391 }
392
393 /* ..... 0 or 欠損値以外の実データの取得 ..... */
394 ax = 0;
395 if (dl)
396 { ax++;
397 if (bh)
398 { bl++; al *= 2; dh = 0;
399 if (bl == 8) { al = *si; si++; rec++; bl = 0; }
400 if (al < 0) dh = 1;
401 }
402 dl--;
403 if (dl)
404 { for (i = 0; i < dl; i++)
405 { bl++; al *= 2; ax *= 2;
406 if (bl == 8) { al = *si; si++; rec++; bl = 0; }
407 if (al < 0) ax++;
408 }
409 }
410 if (dh) ax = -ax;
411 }
412 *di = ax; di++;
413 }
414 *di = d_missing; *bx = 0;
415 return( rec );
416 }

```